

InfoSeg: Unsupervised Semantic Image Segmentation with Mutual Information Maximization

Robert Harb^(✉) and Patrick Knöbelreiter

Institute of Computer Graphics and Vision, Graz University of Technology, Austria
robert.harb@icg.tugraz.at

Abstract. We propose a novel method for unsupervised semantic image segmentation based on **mutual information maximization between local and global high-level image features**. The core idea of our work is to **leverage recent progress in self-supervised image representation learning**. Representation learning methods compute a single high-level feature capturing an entire image. In contrast, we compute multiple high-level features, each capturing image segments of one particular semantic class. To this end, we propose a novel two-step learning procedure comprising a segmentation and a mutual information maximization step. In the first step, we segment images based on local and global features. In the second step, we maximize the mutual information between local features and high-level features of their respective class. For training, we provide solely unlabeled images and start from random network initialization. For quantitative and qualitative evaluation, we use established benchmarks, and COCO-Persons, whereby we introduce the latter in this paper as a challenging novel benchmark. InfoSeg significantly outperforms the current state-of-the-art, e.g., we achieve a relative increase of 26% in the Pixel Accuracy metric on the COCO-Stuff dataset.

Keywords: Unsupervised Semantic Segmentation · Representation Learning.

1 Introduction

Semantic image segmentation is the task of assigning a class label to each pixel of an image. Various applications make use of it, including autonomous driving, augmented reality, or medical imaging. As a result, a lot of research was dedicated to semantic segmentation in the past. However, the vast majority of research focused on supervised methods. A major drawback of supervised methods is that they require large labeled training datasets containing images together with pixel-wise class labels. These datasets have to be created manually by humans with great effort. For example, annotating a single image of the Cityscapes [7] dataset required 90 minutes of human labor on average. This dependence of supervised methods on large human-annotated training datasets limits practical applications. We tackle this problem by introducing a novel approach on semantic image segmentation that does not require any labeled training data.



Fig. 1. (a) Input image. The two magnified image patches have vastly different low-level appearance despite covering the same semantic object: a person. (b) Color based segmentation fails to capture any high-level structure of the image. (c) Representation learning captures high-level information of the entire image in a single feature. (d) InfoSeg captures semantically similar image areas in separate features.

The major challenge of semantic image segmentation is to identify high-level structures in images. State-of-the-art methods approach this by learning from labeled data. While extensive research exists in segmentation without labeled data, it mainly focuses on non-learning based methods using low-level features such as color or edges [11, 6, 18, 1]. In general, low-level features are insufficient for semantic segmentation. They are not homogenous across high-level structures. Figure 1(a-b) illustrate this problem. An image depicting a person is segmented based on color. Color changes vastly across image areas, even if they are semantically correlated. Consequently, the resulting segmentation does not capture any high-level structures. Contrarily, Figure 1(d) illustrates how InfoSeg maps unlabeled images to segmentations that capture high-level structures. These segmentations often directly capture the semantic classes of labeled datasets.

The core idea of our method is to leverage image-level representation learning for pixel-level segmentation. Only recently, self-supervised representation learning methods [14, 5, 26] showed how to extract high-level features from images without any annotated training data. However, they compute features that capture the *entire* content of images. Therefore, they are not suitable for segmentation. To enable segmentation, we instead use multiple high-level features, each capturing semantically similar image areas. This allows us to assign pixels to classes based on their attribution to each of these features. Figure 1(c-d) illustrate how our approach differs from image-level representation learning. We learn high-level features with a mutual information (MI) maximization approach, inspired by Local Deep InfoMax [14]. However, unlike Local Deep InfoMax, we follow a novel two-step learning procedure enabling segmentation. At each iteration, we perform a Segmentation and Mutual Information Maximization step. In the first step, we segment images using the current features. In the second step, we update the features based on the segmentation from the first step. This two-step procedure allows us to train InfoSeg using solely unlabeled images and without pre-trained network backbones.

We motivate the exact structure of InfoSeg by first giving a thorough review of current-state-of-the-art methods [17, 27], followed by a discussion of their limitations and how we approach them in InfoSeg. Our qualitative and quantitative evaluation show that InfoSeg significantly outperforms all compared methods.

虽然有一些分割方法是基于无标注数据的，但他们主要是一些非学习的方法，主要基于low-level feature，比如color或者edges

For example, we achieve a relative increase of 26% in Pixel-Accuracy (PA) on the COCO-Stuff dataset [4]. Even though we follow the standard evaluation protocol for quantitative evaluation, we provide a critical discussion of it and uncover problems left undiscussed by recent work [17,27]. Furthermore, in addition to established datasets, we introduce COCO-Persons as a novel benchmark. COCO-Persons contains complex scenes requiring high-level interpretation for segmentation. Our experiments show that InfoSeg handles the challenging scenes of COCO-Persons significantly better than compared methods. Finally, we perform an ablation study.

2 Related work

Self-Supervised Image Representation Learning aims to capture high-level content of images without using any labeled training data. State-of-the-art methods follow a contrastive learning framework [26,13,14,2,5,10,30,12]. In contrastive learning, one computes multiple representations of differently augmented versions of the same input image. Augmentations can include photometric or geometric image transformations. During training, one enforces similarity on representations computed from the same image and dissimilarity on representations of different images. To this end, various objectives exist, such as the normalized cross entropy [26] or MI [14].

一般的self-supervised image representation learn: 在图片上进行数据增强操作, 训练期间, 强制同一个图片 (及其数据增强的图片) 学到的表达相似, 而不同图片学到的表达不相似

互信息

Unsupervised Semantic Image Segmentation. Invariant Information Clustering (IIC) [17] is a clustering approach also applicable for semantic segmentation. Briefly, IIC uses a MI objective that enforces the same prediction for differently augmented image patches. The authors of IIC proposed to use photometric or geometric image transformations to compute augmentations. For example, one can create augmentations by random color jittering, rotation, or scaling. Ouali *et al.* [27] did a follow-up work on IIC. In addition to standard image transformations, they proposed to process image patches through various masked convolutions. We further discuss these two methods and its differences to InfoSeg in Section 3.2. Concurrent to our work, Mirsadeghi *et al.* proposed InMARS [23]. InMARS is also related to IIC. However, instead of operating on each pixel individually, InMARS utilizes a superpixel representation. Furthermore, a novel adversarial training scheme is introduced.

Another recently introduced method that states to perform unsupervised semantic segmentation is SegSort [15]. However, we note that SegSort still uses supervised learning at multiple stages. First, they initialize parts of their network architecture with pre-trained weights obtained by supervised training of a classifier on the ImageNet [8] dataset. Second, they use pseudo ground truth masks generated by a HED contour detector [31], which is trained supervised using the BSDS500 [1] dataset. Therefore, we do not consider SegSort as an unsupervised method.

3 Motivation

In this section, we first review how recent work [17,27] uses MI for unsupervised semantic image segmentation. Then, we discuss limitations of these methods, and how we tackle them in InfoSeg.

3.1 Unsupervised Semantic Image Segmentation

State-of-the-art methods [17,27] adapt the MI based image clustering approach of IIC [17] for segmentation. In the following, we introduce IICs' approach on image clustering and then the proposed modifications for segmentation.

For clustering, one creates two versions x and x' of the same image. These versions show the same semantic content, but alter low-level appearance by using random photometric or geometric transformations. Consequently, semantic class predictions y and y' of the two images x and x' should be the same. To achieve this, one maximizes the MI between y and y'

$$\max_{\psi} I(\Phi_{\psi}(x); \Phi_{\psi}(x')) = I(y; y'), \quad (1)$$

where Φ is a CNN parametrized by ψ . Considering we can express the MI between y and y' as

$$I(y; y') = H(y) - H(y|y'), \quad (2)$$

Equation (1) maximizes the entropy $H(y)$ while minimizing the conditional entropy $H(y|y')$. Minimizing $H(y|y')$ pushes predictions of the two images x and x' together. Therefore, the network has to compute predictions invariant to the different low-level transformations. This should encourage class predictions to depend on high-level image content instead. While sole minimization of $H(y|y')$ can trivially be done by assigning the same class to all images. Additional maximization of $H(y)$ has a regularization effect against such degenerate solutions. Since maximizing $H(y)$ encourages predictions that put equal probability mass on all classes. Consequently, predictions for all images can not collapse to a single class.

For segmentation, Ji *et al.* [17] proposed to use the previously introduced clustering approach on image patches rather than entire images. Two image versions are pushed through a network that computes dense pixel-wise class predictions. The objective given in Equation (1) is now applied on the pixel-wise class predictions. Therefore, each prediction depends on an image patch rather than an entire image. Patches are defined by the receptive field for each output pixel of the network. Additionally, one enforces local spatial invariance by maximizing MI of predictions from adjacent image patches. This approach on unsupervised semantic segmentation was initially proposed by IIC [17]. Furthermore, Ouali *et al.* [27] proposed an extension by generating views using different masked convolutions [25]. In the following, we discuss three major limitations of these two works, and how we tackle them in InfoSeg.

首先对同一张图片生成两个版本，由于两个版本包含的语义内容是一致的，因此输出的semantic prediction也应该是一致的

网络最大化 y 和 y' 之间的MI，MI可以写为 y 的熵减去 y' 条件下的 y 的条件熵，后者使得 y 和 y' 的预测结果尽量相近，但单独优化该项会引起退化：即所有的图像都预测为一个结果即可，前者可以看做后者的一个正则项，即使得所有的 y 能尽可能的平均分布在所有类别上

对于分割任务，其实是把分类任务的image level换成了pixel level

3.2 Limitations of current methods

The first limitation of discussed methods is that **they do not incorporate global image context**. Global context is essential to capture high-level structures, since they often cover large image areas having diverse local appearance. Therefore observing only small image patches is often not sufficient to identify them. Ideally, each pixel-wise prediction should depend on the entire image. Nevertheless, the discussed approaches make pixel-wise predictions based on image patches. The receptive field of the network Φ determines the size of these patches. In general, one could enlarge the receptive field by changing the network architecture. However, adapting IIC from clustering to segmentation is based on restricting each prediction’s receptive field from entire images to patches. By making each pixel-wise prediction dependent on the entire image again, one would fall back to clustering. In InfoSeg we capture global context in global high-level features that cover the entire image. We make pixel-wise predictions based on the MI between these global features and local patch-wise features. This allows each pixel-wise prediction to depend on the entire image.

A second limitation of discussed methods is that they **fail to leverage recent advances in image representation learning** [14,5,26,2]. These methods are effective at capturing high-level image content, but only at the image-level. Adapting them for pixel-level segmentation is not trivial. Ouali *et al.* [27] attempted this with their Autoregressive Representation Learning (ARL) loss, but failed to increase segmentation performance. Despite high-level information is constant across large image areas, ARL computes for *each pixel* a separate high-level feature. Contrarily, in InfoSeg, we share high-level features over the *entire image*. To still allow pixel-wise segmentation, we compute multiple high-level features. Each high-level feature encodes only image areas depicting one class. We then assign pixels to classes based on their attribution to each of these features.

Finally, **discussed methods jointly learn features and segmentations**. They use intermediate feature representations to assign pixels to class labels. At the beginning of training, features depend on random initialization and contain no high-level information. This can lead to classes that latch onto low-level features instead of capturing high-level information. This issue was first discussed for image classification by SCAN [29]. Instead, **we decouple feature learning and segmentation**. Therefore, we perform two steps at each iteration. First, we compute features that are explicitly trained to encode high-level information. Then, we use them for segmentation.

4 InfoSeg

In InfoSeg, we tackle unsupervised semantic image segmentation. We take a set $\{\mathbf{X}^{(n)} \in \mathcal{X}\}_{n=1}^N$ of N unlabeled images and assign a label $\mathcal{Z} = \{z_1, \dots, z_K\}$ to every pixel of each image. Importantly, for one particular image, we do not specify which nor how many labels should be assigned. **We only provide the total number of labels K in all images**. After training, we follow the standard

对于每张图片，并不指定该图片中有多少种语义label，而只给定数据集中的label类别

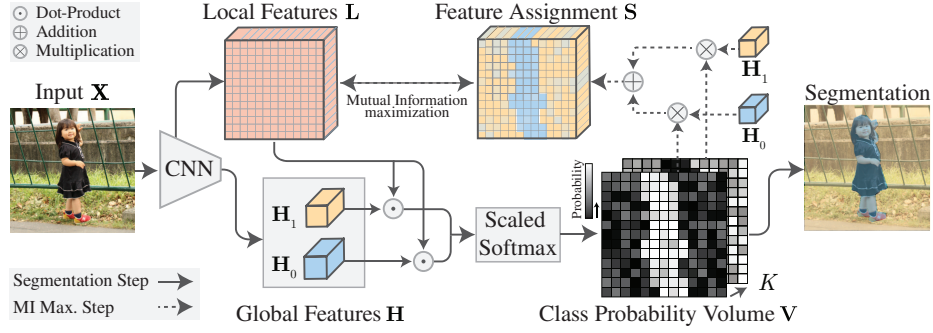


Fig. 2. Overview of InfoSeg for $K = 2$ classes. At each training iteration, we alternate the following two steps. **Segmentation Step** (solid lines): An input image \mathbf{X} is passed through a CNN to compute local patch-wise features \mathbf{L} and for each class k a global image-level feature \mathbf{H}_k . We then score local \mathbf{L} with global \mathbf{H} features using a dot-product. The result is passed through a scaled softmax function to compute the class probability volume \mathbf{V} . Finally, we obtain a segmentation by assigning each pixel to the class with the largest probability. **Mutual Information Maximization Step** (dashed lines): The global feature assignment \mathbf{S} is computed as a sum of global features, weighted by their respective class probabilities at each spatial position. Finally, we maximize Mutual Information between local features \mathbf{L} and their respective feature assignment \mathbf{S} .

训练结束后，将学到的label映射到人工标注的label上

evaluation protocol and map the learned labels of InfoSeg directly to the semantic classes of an annotated dataset.

InfoSeg is designed to tackle the three limitations of state-of-the-art methods discussed in Section 3.2. Figure 2 shows an overview of InfoSeg. In the following, we first discuss how we leverage recent progress in representation learning for semantic segmentation in Section 4.1. Then we provide further details of our method in Section 4.2 and Section 4.3.

4.1 Representation Learning for Segmentation

We first review how Local Deep InfoMax [14] captures high-level information of entire images, and then how InfoSeg adapts this approach to target image segmentation.

Local Deep InfoMax [14] learns global high-level features of images by maximizing their average MI with local features. Local features cover image patches, and the global feature covers the entire image. If the global feature has limited capacity, the network cannot simply copy all local features' content into the global feature to maximize MI. Instead, the network has to encode a compact representation that shares information with as many image patches as possible. Hjelm *et al.* [14] showed that the resulting global features encode high-level image information. They motivated this by the idea that high-level information is

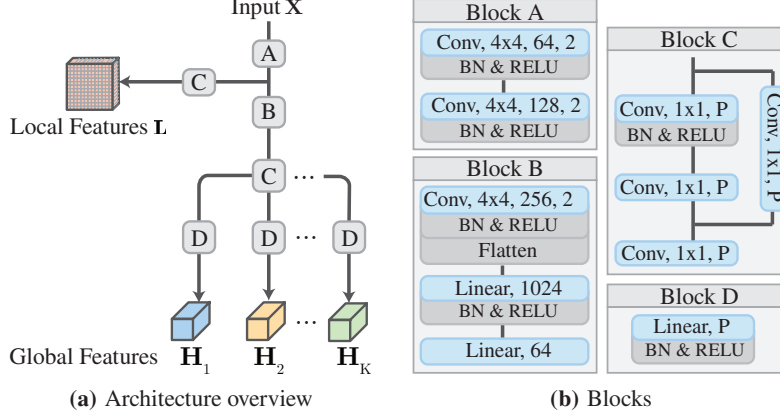


Fig. 3. Feature computation for K classes. (a) Overview of network architecture. (b) Used blocks. **Legend:** *Conv*, $W \times W, C, d$: Convolution with filter size $W \times W$, C channels and stride d . Blocks that are used multiple times, each have their own set of parameters.

often constant over an entire image, while low-level information such as pixel-level noise varies. Consequently, the global feature is encouraged to encode the former while disregarding the latter.

To enable pixel-wise segmentation, we compute for each image multiple global features instead of a single one. Each global feature only encodes image areas that depict a particular class. This allows us to segment images by assigning pixels to classes based on their attribution to each global feature. During training, we maximize for each global feature MI only with local features covering its respective class. Therefore, we learn high-level features in a similar way as Local Deep InfoMax [14], but target segments instead of entire images. This requires us to learn high-level features together with segmentations. To this end, we alternate two steps at each iteration. In the *Segmentation Step*, we assign local to global features based on their content, *i.e.*, we segment images. In the *Mutual Information Maximization step*, we maximize the MI between all global features and assigned local features, *i.e.*, we learn the features. We describe both steps in the following.

4.2 Segmentation Step

Given an input image $\mathbf{X} \in \mathcal{X} = \mathbb{R}^{M \times N \times C}$, we compute P -dimensional global $\mathbf{H} \in \mathbb{R}^{K \times P}$ and patch-wise local $\mathbf{L} \in \mathbb{R}^{U \times V \times P}$ features. The k -th global feature $\mathbf{H}_k \in \mathbb{R}^P$ encodes a high-level representation for the k -th class and covers the entire image. The local feature $\mathbf{L}_{i,j} \in \mathbb{R}^P$ at the spatial position (i, j) encodes an image patch. Furthermore, the spatial resolution of \mathbf{L} is downsampled by a rate of d from the input resolution, *i.e.* $U = M/d$ and $V = N/d$.

Figure 3 shows the architecture of our feature computation network. First, the input image is processed by Block A, resulting in a grid of patch-wise image features. To compute the local features \mathbf{L} , we further process these patch-wise features by Block C. Adding this additional residual block of pointwise convolutions led to better performance, than using Block A’s output directly for the local features. To compute the global features \mathbf{H}_k , we first process the output of Block A to image-level features using Block B. Then, similarly, as for the local features, we add a residual block of pointwise convolutions using Block C. Finally, each global feature is computed using a separate linear layer using Block D.

To compute an image segmentation, we use the dot-product of a local and global feature pair $\langle \mathbf{L}_{i,j}, \mathbf{H}_k \rangle$ as a class score. A high score indicates that the k -th class is shown at the position (i, j) . We elaborate in Section 4.3 how MI maximization increases the dot-product of a local feature and the global feature of its corresponding class. After computing the class scores, we apply a pixel-wise scaled softmax to compute a class-probability volume $\mathbf{V} \in \mathbb{R}^{U \times V \times K}$ with elements

$$V_{i,j,k} = \frac{\exp(\tau \cdot \langle \mathbf{L}_{i,j}, \mathbf{H}_k \rangle)}{\sum_k \exp(\tau \cdot \langle \mathbf{L}_{i,j}, \mathbf{H}_k \rangle)}, \quad (3)$$

where τ is a hyper-parameter that controls the smoothness of the resulting distribution. Using the probability volume, we compute the low-res segmentation \mathbf{K} for every pixel (i, j) with

$$k_{i,j} = \arg \max_{k \in \mathcal{Z}} V_{i,j,k}, \quad (4)$$

by taking the class with the largest probability. We can then compute a full-res segmentation \mathbf{Z} by upsampling the low-res segmentation \mathbf{K} to the input image resolution.

4.3 Mutual Information Maximization Step

We first need to assign each local feature to its corresponding class’s global feature. We could do this using the segmentation \mathbf{K} . However, this disregards class probabilities, instead of utilizing their exact values, *e.g.* to account for uncertainty. Especially at the beginning of training, segmentations are uncertain and depend on random network initialization. Reinforcing possibly incorrect predictions can lead to degenerate solutions. To alleviate this problem, we do not make hard class assignments using \mathbf{K} , but soft assignments using class-probabilities \mathbf{V} . Instead of assigning a single global feature, we weight each global feature by its respective class probability. To this end, we define the function $S_{\theta}^{(i,j)}$ that computes a soft global feature assignment for the local feature $\mathbf{L}_{(i,j)}$ as follows

$$S_{\theta}^{(i,j)}(\mathbf{X}) = \sum_k V_{i,j,k} \cdot H_{\theta}^{(k)}(\mathbf{X}), \quad (5)$$

where the function $H_{\theta}^{(k)}(\mathbf{X})$ computes the k -th global feature \mathbf{H}_k for an image \mathbf{X} , and θ denotes the learnable parameters of our network.

对每个 local feature，计算与每个 class 的 global feature 的点乘，（共 K 个 class）并在 K 维度上进行 softmax，得分最高的就定为该类别

（因为是在 \mathbf{L} 上做的，而 \mathbf{L} 是由输入图片卷积得到的 patch feature，因此分辨率较低，也就是这一步得到了低分辨率的分割图像）

之后，全分辨率的分割结果可以通过上采样低分辨率分割结果得到

segmentation \mathbf{K} 隐去了每个像素点属于每个 class 的概率，特别是在训练初期，segmentation 是取决于随机初始化的，可能会引起退化的结果，因此，这里并没有使用 \mathbf{K} ，而是使用 \mathbf{V} ，也就是 softmax 后的得到的概率图

对于每个像素位置，使用预测到的属于每一类的概率值加权每一类的 global feature

During training, we maximize the MI between the output of $S_{\theta}^{(i,j)}(\mathbf{X})$ and the corresponding local feature $\mathbf{L}_{(i,j)}$ for all spatial positions (i, j) . Hence our objective is given as

$$\max_{\theta} \mathbb{E}_{\mathbf{X}} \left[\frac{1}{UV} \sum_{i,j} I \left(L_{\theta}^{(i,j)}(\mathbf{X}); S_{\theta}^{(i,j)}(\mathbf{X}) \right) \right], \quad (6)$$

where $\mathbb{E}_{\mathbf{X}}$ denotes the expectation over all training images \mathbf{X} and the function $L_{\theta}^{(i,j)}(\mathbf{X})$ computes the local features $\mathbf{L}_{i,j}$ given an input image \mathbf{X} .

To evaluate our objective Equation (6), consider that local and global features are high-dimensional continuous random variables. MI computation of such variables is challenging. Contrarily to discrete variables as in the objective of IIC Equation (1), where exact computation is possible. For continuous variables, Belghazi *et al.* [3] proposed MI estimation by maximizing lower bounds parametrized by neural networks. They used a bound based on the Donsker & Varadhan (DV) representation of the Kullback-Leibler (KL) divergence. While several other bounds exist [28], we use a bound based on the Jensen-Shannon Divergence (JSD). Mainly because Hjelm *et al.* [14] showed favorable properties of the JSD bound compared to others in their representation learning setting. This includes increased training stability and better performance with smaller batch sizes. Nevertheless, we also perform experiments using the DV bound in our ablation studies. A JSD based MI estimator $\hat{I}_{JSD}(X; Y)$ for two random variables X and Y can be defined as follows [24]

$$I(X; Y) \geq \hat{I}_{JSD}(X; Y) := \mathbb{E}_{p(x,y)}[-\text{sp}(-T(x, y))] - \mathbb{E}_{p(x)p(y)}[\text{sp}(T(x, y))], \quad (7)$$

where $\text{sp}(x) = \log(1 + e^x)$ and T is a discriminator mapping sample pairs from X and Y to a real valued score. The first and second expectations are taken over samples from the joint $p(x, y)$ and marginal $p(x)p(y)$ distributions. Consequently, to tighten the bound, the discriminator T needs to discriminate samples from the joint and marginal distributions by assigning high or low scores, respectively.

To use the JSD estimator Equation (7) in our objective Equation (6), we have to define the discriminator T and a sampling strategy. Following recent work [14, 2], we create joint and marginal samples by combining feature pairs computed from the same image \mathbf{X} and two randomly paired images \mathbf{X} and \mathbf{X}' , respectively. The discriminator T can be implemented using any arbitrary function that maps feature pairs to a discrimination score, *e.g.*, a neural network. For efficiency, we use the dot-product to compute discrimination scores *i.e.*, $T(x, y) := \langle x, y \rangle$. This requires only a single expensive forward pass through our network to compute the features, while we can then score any arbitrary combination with a cheap dot-product. Omitting the spatial indices (i, j) to avoid notational clutter, this leads to the MI estimator

$$\hat{I}_{JSD}(L_{\theta}(\mathbf{X}); S_{\theta}(\mathbf{X})) := \mathbb{E}_{\mathbb{P}}[-\text{sp}(-\langle L_{\theta}(x), S_{\theta}(x) \rangle)] - \mathbb{E}_{\mathbb{P} \times \mathbb{P}}[\text{sp}(\langle L_{\theta}(x); S_{\theta}(x') \rangle)], \quad (8)$$

where \mathbb{P} is the empirical distribution of our dataset, x is an image sampled from \mathbb{P} and x' is an image sampled from $\hat{\mathbb{P}} = \mathbb{P}$. We can now simply insert the estimator of Equation (8) into our objective Equation (6). Maximizing the resulting objective increases the dot-product of local features with the global feature of their assigned class. Consequently, we use the dot-product as a class score, as described in Section 4.2.

5 Experiments

We first introduce our experimental setup and discuss challenges at the quantitative evaluation of unsupervised segmentation. Then we perform an evaluation using established benchmarks [22,4], and COCO-Persons, a novel dataset introduced in this work. On all datasets, InfoSeg significantly outperforms compared methods. Finally, we perform ablation studies.

5.1 Setup

We start training from random network initialization and provide solely unlabeled images. We set $P = 1024$, $\tau = 0.8$ and use the ADAM optimizer [19] with a learning rate of 10^{-4} and a batch size of 64. Furthermore, the network architecture we use results in a downsampling rate of $d = 4$, and we set the number of classes K to be equal to the number of classes in each dataset.

Note that InfoSeg requires a network with a different structure as Invariant Information Clustering (IIC) and Autoregressive Clustering (AC). For InfoSeg, the final outputs are 1×1 sized global image features. Contrarily, in IIC and AC, the final outputs are pixel-wise class predictions downsampled from the input image resolution. This impedes a comparison with these methods using the exact same architecture. Nevertheless, we provide an experiment in our ablation study where we apply the objective of IIC on the output of our Segmentation Step.

5.2 Quantitative Evaluation

Meaningful quantitative evaluation of unsupervised semantic segmentation is challenging. Recent work used the PA for quantitative evaluation. The PA is defined as the percentage of pixels assigned to the same class as in a given annotation. However, in unsupervised semantic segmentation, one does not specify which classes should be used for segmentation. Instead, many different segmentations can be considered as equally valid. Nevertheless, quantitative evaluation metrics, such as the PA or mean Intersection-Over-Union (mIoU), evaluate all pixel-wise predictions as incorrect that do not exactly match the given annotations. While this has been left undiscussed by previous work [17,27], we emphasize this has to be considered when interpreting quantitative metrics of unsupervised methods.

We can further illustrate problems at quantitative evaluation using the COCO-Stuff [4] dataset as an example. The dataset contains the class *rawmaterial* that

labels image areas depicting metal, plastic, paper, or cardboard. We argue that this is a very specific class and aggregating these four materials in one class is an arbitrary design choice of the dataset. It is unfeasible to expect an unsupervised method to come up with this specific solution. Nevertheless, recent methods [17,27] reported significant increases over baseline models on the PA. We attribute this to the dataset’s vast class imbalance. Besides very specific classes such as *rawmaterial*, the dataset also contains more generic classes such as *water*, or *plant*. These classes are overrepresented and make up more than 50% of all pixels. Therefore, **an algorithm can achieve high PA by focusing mainly on these few overrepresented classes.** To illustrate this effect, we provide a confusion matrix of our predictions in the supplementary material.

Despite the discussed problems, we follow prior work and use the PA to evaluate all of our results quantitatively. Following the standard evaluation protocol [17,27], **we map each of the predicted classes in \mathcal{Z} to one of the annotated classes in \mathcal{Z}' before computing the PA.** This is necessary because class ordering is unknown without providing labeled data during training. **We find the one-to-one mapping between \mathcal{Z} and \mathcal{Z}' by solving the linear assignment problem using the Hungarian method [20].** We compute this mapping once after training and use the same mapping for all images in the dataset.

首先将预测类别中的每个类映射到标注类别中的一个类，这种映射是通过使用 Hungarian 方法解决 linear assignment 问题实现的

5.3 Data

Recent work [17,27] established the COCO-Stuff [4] and Potsdam [22] datasets as benchmarks. COCO-Stuff contains 15 classes and Potsdam 6 classes. Additionally, for both datasets, a reduced 3-class variation exists. We use the same pre-processing as in the compared methods, resulting in 128×128 sized RGB images for COCO-Stuff, and 200×200 sized RGBIR images for Potsdam.

While unsupervised segmentation of COCO-Stuff and Potsdam is challenging, most classes in these datasets still have a homogeneous low-level appearance. For example, low-level features such as color and texture are often sufficient to segment areas labeled as *water* in COCO-Stuff or *road* in Potsdam. To show that InfoSeg can go one step further, we evaluate on an additional dataset where segmentation is more reliant on high-level image features. To this end, we introduce the COCO-Persons dataset, which we will provide publicly. Each image depicts one or multiple persons and is annotated with a person and a non-person class. Face, hair, and clothing of persons vary vastly in color, texture, and shape, and the non-person areas cover a variety of complex indoor and outdoor scenes. The dataset is a subset of the COCO [21] dataset and contains 15 399 images having 128×128 pixels.

COCO-Stuff 和 Potsdam 中的大部分类别都可以通过 low-level 的 appearance 来识别，比如 water 和 road，可以直接通过颜色和纹理识别出来，因此这里引入了一个新的更具有挑战性的数据集，COCO-Persons

5.4 Results

We provide quantitative and qualitative results in Table 1 and Figure 4, respectively. To compute results for COCO-Persons, we used publicly available

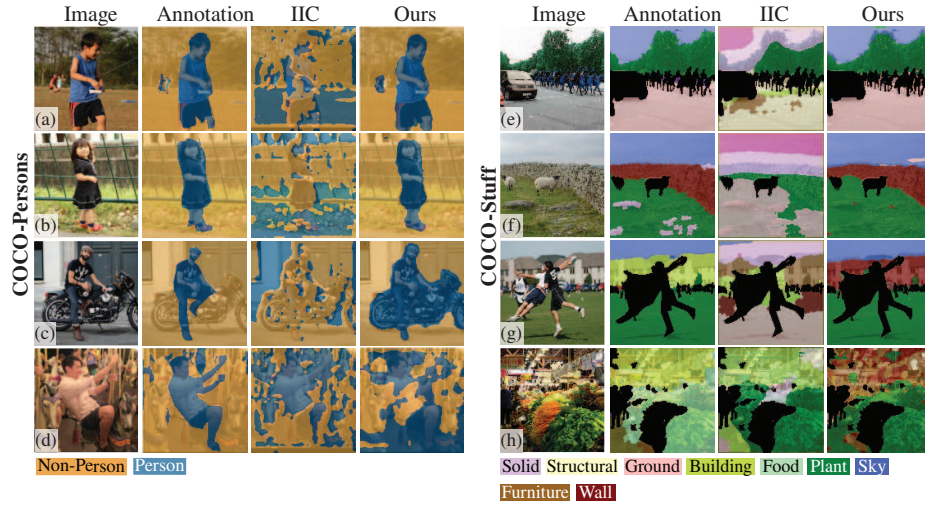


Fig. 4. Qualitative comparison. Non-stuff areas in COCO-Stuff are masked in black.

Method	COCO-Persons	COCO-Stuff	COCO-Stuff-3	Potsdam	Potsdam-3
Random CNN	52.3	19.4	37.3	28.3	38.2
K-Means	54.3	14.1	52.2	35.3	45.7
Doersch* [9]	55.6	23.1	47.5	37.2	49.6
Isola* [16]	57.5	24.3	54.0	44.9	63.9
IIC [17]	57.1	27.7	72.3	45.4	65.1
AC [27]	-	30.8	72.9	49.3	66.5
InMARS [23]	-	31.0	73.1	47.3	70.1
InfoSeg (ours)	69.6	38.8	73.8	57.3	71.6

Table 1. Pixel-Accuracy of InfoSeg and compared methods. *Clustering of features from methods that are not specifically designed for image segmentation.

implementations, if available. In our experiments, InfoSeg significantly outperformed all compared methods [9,16,17,27,23]. We discuss qualitative results in the following.

COCO-Persons. In Figure 4(a-c), we show successful segmentation of images with vastly inhomogenous low-level appearance. InfoSeg even captures the two small persons in the background of Figure 4(a). In Figure 4(c), the motorbike is assigned to the same class as the person. The dataset contains several images where persons are shown together with motorbikes. Therefore, without supervision, it is challenging to disentangle these two semantic concepts. In Figure 4(d), we show a challenging example yielding a failure case.

COCO-Stuff. Figure 4(e-f) show examples where our predictions are close to the annotations. Figure 4(g-h) provide reasonable segmentations, even though large

portions differ from the annotations. These examples demonstrate challenges at the evaluation of COCO-Stuff due to overly specific classes. Matching the annotations requires precise distinction of similar high-level concepts, which is difficult without supervision. The example in Figure 4(g) shows multiple houses that are assigned to the same class as the stone wall in Figure 4(f). However, the ground truth of COCO-Stuff assigns the stone wall to a *wall* class and the houses to a *building* class. Figure 4(h) shows a market scene containing vegetables labeled as *food* but predicted as *plants*. Arguably, vegetables are food and plants.

5.5 Ablation Studies

To examine the influence of individual components, we perform the following ablation studies. First, we evaluate the effectiveness of soft assignments **by replacing them with hard assignments**. Therefore, we change our objective Equation (6) to maximize the MI at each spatial position between the local feature and the global feature of the assigned class according to the segmentation **K**. Second, we **replace the JSD MI estimator with a DV one**. Finally, in the last ablation

study, we omit our Mutual Information Maximization step and solely perform our Segmentation Step. As a replacement for our Mutual Information Maximization step we apply the MI maximization objective of IIC, referred to as IIC-MI. To create the two image versions required by IIC-MI, we use the same transformations as in IIC.

Table 2 shows the results of our ablation studies, whereby we performed all experiments using the COCO-Stuff-3 dataset. We can observe the following: Using soft assignments increases performance over hard assignments. A JSD-based MI estimator performs better than a DV-based, which aligns with the results of Hjelm *et al.* [14]. And replacing our Mutual Information Maximization step with the objective of IIC leads to a decline in performance.

	Measure	MI Max. Step	PA
IIC-MI	-		60.2
DV	Hard Assignment		53.9
DV	Soft Assignment		55.1
JSD	Hard Assignment		67.3
JSD	Soft Assignment		73.8

Table 2. Ablation studies on COCO-Stuff-3.

6 Conclusion

We proposed a novel approach for unsupervised semantic image segmentation. Our experiments showed that our method yields semantically meaningful predictions and significantly outperforms related methods. We used the established datasets for evaluation and introduced a novel challenging benchmark COCO-Person. Furthermore, we discussed several problems making the quantitative evaluation of unsupervised semantic segmentation challenging. Finally, we performed ablation studies on our model.

Acknowledgements. This work was partly funded by the Austrian Research Promotion Agency (FFG) under project 874065.

References

1. Arbelaez, P., Maire, M., Fowlkes, C., Malik, J.: Contour detection and hierarchical image segmentation. vol. 33, pp. 898–916. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2010) [2](#), [3](#)
2. Bachman, P., Hjelm, R.D., Buchwalter, W.: Learning representations by maximizing mutual information across views. In: *Conference on Neural Information Processing Systems*. pp. 15535–15545 [3](#), [5](#), [9](#)
3. Belghazi, M.I., Baratin, A., Rajeshwar, S., Ozair, S., Bengio, Y., Hjelm, D., Courville, A.: Mutual information neural estimation. In: *International Conference on Learning Representations*. pp. 530–539 (2018) [9](#)
4. Caesar, H., Uijlings, J., Ferrari, V.: Coco-stuff: Thing and stuff classes in context. In: *Computer Vision and Pattern Recognition*. pp. 1209–1218 (2018) [3](#), [10](#), [11](#)
5. Chen, T., Kornblith, S., Norouzi, M., Hinton, G.: A simple framework for contrastive learning of visual representations. In: *arXiv preprint arXiv:2002.05709* (2020) [2](#), [3](#), [5](#)
6. Comaniciu, D., Meer, P.: Mean shift: A robust approach toward feature space analysis. pp. 603–619. No. 5, *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2002) [2](#)
7. Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., Schiele, B., R&d, D.A., Darmstadt, T.U.: The Cityscapes Dataset for Semantic Urban Scene Understanding. Tech. rep., www.cityscapes-dataset.net [1](#)
8. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: ImageNet: A Large-Scale Hierarchical Image Database. In: *Computer Vision and Pattern Recognition* (2009) [3](#)
9. Doersch, C., Gupta, A., Efros, A.A.: Unsupervised visual representation learning by context prediction. In: *International Conference on Computer Vision*. pp. 1422–1430 (2015) [12](#)
10. Dosovitskiy, A., Springenberg, J.T., Riedmiller, M., Brox, T.: Discriminative unsupervised feature learning with convolutional neural networks. In: *Conference on Neural Information Processing Systems*. pp. 766–774 (2014) [3](#)
11. Felzenszwalb, P.F., Huttenlocher, D.P.: Efficient graph-based image segmentation. vol. 59, pp. 167–181. *International Journal of Computer Vision* (2004) [2](#)
12. He, K., Fan, H., Wu, Y., Xie, S., Girshick, R.: Momentum contrast for unsupervised visual representation learning. In: *Computer Vision and Pattern Recognition* (2020) [3](#)
13. Henaff, O.: Data-efficient image recognition with contrastive predictive coding. In: *International Conference on Machine Learning*. pp. 4182–4192 (2020) [3](#)
14. Hjelm, R.D., Fedorov, A., Lavoie-Marchildon, S., Grewal, K., Bachman, P., Trischler, A., Bengio, Y.: Learning deep representations by mutual information estimation and maximization. In: *International Conference on Learning Representations* (2019) [2](#), [3](#), [5](#), [6](#), [7](#), [9](#), [13](#)
15. Hwang, J.J., Yu, S.X., Shi, J., Collins, M.D., Yang, T.J., Zhang, X., Chen, L.C.: Segsort: Segmentation by discriminative sorting of segments. In: *International Conference on Computer Vision* (2019) [3](#)
16. Isola, P., Zoran, D., Krishnan, D., Adelson, E.H.: Learning visual groups from co-occurrences in space and time. In: *arXiv preprint arXiv:1511.06811* (2015) [12](#)
17. Ji, X., Henriques, J.F., Vedaldi, A.: Invariant information clustering for unsupervised image classification and segmentation. In: *International Conference on Computer Vision* (October 2019) [2](#), [3](#), [4](#), [10](#), [11](#), [12](#)

18. Jianbo Shi, Malik, J.: Normalized cuts and image segmentation. In: IEEE Transactions on Pattern Analysis and Machine Intelligence. vol. 22, pp. 888–905 (Aug 2000) [2](#)
19. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: International Conference on Learning Representations (2015) [10](#)
20. Kuhn, H.W.: The hungarian method for the assignment problem. In: Naval research logistics quarterly. vol. 2, pp. 83–97. Wiley Online Library (1955) [11](#)
21. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: European Conference on Computer Vision. pp. 740–755 (2014) [11](#)
22. Markus Gerke, I.: Use of the stair vision library within the isprs 2d semantic labeling benchmark (vaihingen) [10](#), [11](#)
23. Mirsadeghi, S.E., Royat, A., Rezaatofghi, H.: Unsupervised image segmentation by mutual information maximization and adversarial regularization. IEEE Robotics and Automation Letters volume 6 (4), 6931–6938 (2021) [3](#), [12](#)
24. Nowozin, S., Cseke, B., Tomioka, R.: f-gan: Training generative neural samplers using variational divergence minimization. In: Conference on Neural Information Processing Systems. pp. 271–279 (2016) [9](#)
25. Van den Oord, A., Kalchbrenner, N., Espeholt, L., Vinyals, O., Graves, A., et al.: Conditional image generation with pixelcnn decoders. In: Conference on Neural Information Processing Systems. pp. 4790–4798 (2016) [4](#)
26. den Oord, A.V., Li, Y., Vinyals, O.: Representation learning with contrastive predictive coding. In: arXiv preprint arXiv:1807.03748. vol. abs/1807.03748 (2018) [2](#), [3](#), [5](#)
27. Ouali, Y., Hudelot, C., Tami, M.: Autoregressive unsupervised image segmentation. In: European Conference on Computer Vision (August 2020) [2](#), [3](#), [4](#), [5](#), [10](#), [11](#), [12](#)
28. Poole, B., Ozair, S., van den Oord, A., Alemi, A., Tucker, G.: On variational bounds of mutual information. In: International Conference on Machine Learning. pp. 5171–5180 (2019) [9](#)
29. Van Gansbeke, W., Vandenhende, S., Georgoulis, S., Proesmans, M., Van Gool, L.: Scan: Learning to classify images without labels. In: European Conference on Computer Vision (2020) [5](#)
30. Wu, Z., Xiong, Y., Yu, S.X., Lin, D.: Unsupervised feature learning via non-parametric instance discrimination. In: IEEE Transactions on Pattern Analysis and Machine Intelligence. pp. 3733–3742 (2018) [3](#)
31. Xie, S., Tu, Z.: Holistically-nested edge detection. In: International Conference on Computer Vision. vol. 125, pp. 3–18. Kluwer Academic Publishers, Hingham, MA, USA (Dec 2017) [3](#)